

First principles view on chemical compound space: Gaining rigorous atomistic control of molecular properties

O. Anatole von Lilienfeld^{1,*}

¹Argonne Leadership Computing Facility, Argonne National Laboratory, Argonne, Illinois 60439, USA

(Dated: October 23, 2012)

A well-defined notion of chemical compound space (CCS) is essential for gaining rigorous control of properties through variation of elemental composition and atomic configurations. Here, we review an atomistic first principles perspective on CCS. First, CCS is discussed in terms of variational nuclear charges in the context of conceptual density functional and molecular grand-canonical ensemble theory. Thereafter, we revisit the notion of compound pairs, related to each other via “alchemical” interpolations involving fractional nuclear chargens in the electronic Hamiltonian. We address Taylor expansions in CCS, property non-linearity, improved predictions using reference compound pairs, and the ounce-of-gold prize challenge to linearize CCS. Finally, we turn to machine learning of analytical structure property relationships in CCS. These relationships correspond to inferred, rather than derived through variational principle, solutions of the electronic Schrödinger equation.

I. INTRODUCTION

In analogy to the vastness and sparseness of outer space, we can loosely refer to the space of chemical systems as chemical compound space (CCS), i.e. some continuous observable space that is populated by all experimentally *and* theoretically possible chemicals with integer nuclear charges and interatomic distances for which chemical interactions occur¹. Stated more precisely, CCS refers to the combinatorial set of all compounds that can be isolated and constructed from possible combinations and configurations of N_I atoms and N_e electrons in real space. In absence of external fields and given N_e and N_I atom-types $\{Z_I\}$ and spatial configurations $\{\mathbf{R}_I\}$, not only covalent, ionic, and metallic bonding result, but also the much weaker hydrogen and van-der-Waals-bonding, responsible for the physics and chemistry of molecular crystals, liquids, and other supra-molecular aggregates, can be derived, as well as all other quantum and statistical mechanical properties such as electronic states, electronic and vibrational spectra, free energies, and even phase diagrams and rare events such as chemical reactions. While most research efforts in this first principles context have been dedicated to the approximations and methods necessary for making property predictions for given compounds, the focus of this tutorial is a first principles view on the compounds *per se*.

Notwithstanding chemical bonding or conformations and merely considering the number of possible stoichiometries it is obvious that the size of CCS is unfathomably large for all but the smallest systems. Due to all the possible combinations of assembling many and various atoms its size scales exponentially with compound size as $\propto Z_{max}^{N_I}$. Here Z_{max} is the number of possible atom types, i.e. the maximal permissible nuclear charge in Mendeleev’s table ($Z_{max} > 100$), and N_I depends on the employed definition of “isolated system” but can certainly reach Avogadro’s number scale for living organisms, chunks of unordered matter, or planets. While many of such speculative compounds are

likely to be unstable, the state of affairs worsens when accounting for the additional degrees of freedom which arise from distinguishable geometries due to differences in atom bonding or conformations. This combinatorial explosion with system size is the main motivation for advocating an *ab initio*, or first principles, view on CCS, i.e. a view that restricts us to use solely $\{Z_I\}$ and $\{\mathbf{R}_I\}$ as input variables¹³⁸, and, while maybe not free of parameters, will not change in its parameterization as $\{Z_I\}$ and $\{\mathbf{R}_I\}$ are freely varied². A major part of modern electronic structure theory and interatomic potential work is concerned with the development of improved methods and approximations for solving Schrödinger’s equation (SE) within the Born-Oppenheimer approximation for Hamiltonians relevant to materials, biological, or chemical research, and deriving properties thereof³. *Ab initio* statistical mechanics efforts are dedicated to sampling the corresponding $3N_I - 6$ degrees of freedom from first principles⁴. In the context of CCS, the electronic Hamiltonian H for solving SE, $H\Psi = E\Psi$, of *any* compound with a given charge, $Q = N_p - N_e$, is uniquely determined by its (unperturbed) external potential, $v(\mathbf{r}) = \sum_I Z_I/|\mathbf{r} - \mathbf{R}_I|$, i.e. by its set $\{\mathbf{R}_I, Z_I\}$. Here, N_p is the total number of protons in the system, i.e. the sum over all nuclear charges. Due to the Hohenberg-Kohn theorem we also know that the electron density $n(\mathbf{r})$, and all electronic properties derived thereof, are determined by $\{Z_I, \mathbf{R}_I\}$, up to a trivial constant, $\{H(\mathbf{r}), N_e\} \leftrightarrow \{Z_I, \mathbf{R}_I, N_e\} \leftrightarrow \{v(\mathbf{r}), N_e\} \leftrightarrow n(\mathbf{r})$ ⁵. Consequently, we work with $\{Z_I, \mathbf{R}_I, N_e\}$.

In this tutorial, CCS is first briefly illustrated in terms of a rough energy scale in section II. In section III we will review the notion of a molecular grand-canonical ensemble density functional theory that can deal with fractional electrons *and* nuclear charges. Section IV will deal with pairs of chemical compounds, and with efforts to exploit the arbitrariness of interpolating functions. It also details the challenge associated to a prize award of one ounce of gold. Finally, we will study recent efforts to use intelligent data analysis methods (machine learning)

to systematically infer analytical structure property relationships from previously calculated electronic structure data sets in section V.

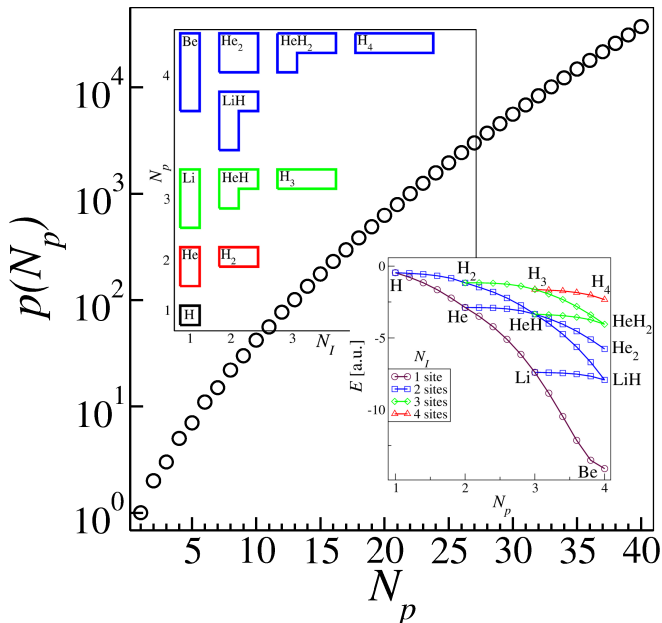


FIG. 1: Exponential scaling of the total number of all possible N partitions, i.e. stoichiometries, as a function of N_p (Regime (i) in section II). Inset upper left-hand side: Young-Ferrers diagrams illustrating the possible partitions (stoichiometries) for N_p as a function of number of atoms, N_I . The color code corresponds to the total number of protons in the compound, $N_p \in \{1 \text{ (black)}, 2 \text{ (red)}, 3 \text{ (green)}, 4 \text{ (blue)}\}$. Inset lower right-hand side: Total potential energy of relaxed molecules as a function of N_p using interpolated pseudopotentials in analogy to Fig. 2 (all systems neutral, BLYP DFT level of theory^{6,7}, arbitrary energy origin due to use of pseudopotentials).

II. ENERGY HIERARCHY

Regarding CCS it is useful to think of a variable system that is comparable to the Mendeleev’s table of the elements. Compounds, however, have many more dimensions than a single atom’s nuclear charge, specifically $4N_I - 6$ ($4N_I - 5$ if linear). One way of thinking about CCS is in terms of an abstract *Gedankenexperiment* involving all theoretically existing compounds. Consider all the compounds possible for a set of protons, subject to a varying amount of kinetic energy (or temperature), provided by a thermostat. Regimes emerge of various familiar degrees of freedom. In this “phase diagram” of CCS these various regimes correspond to

- (i) *stoichiometrical isomers*: A fictitious “very high” temperature regime. Let us assume such high temperatures that all bonds break, and that all spatial degrees of freedom can safely be neglected. Fur-

thermore, we assume isomers to have the same number of elementary particles, N_p and N_e . How many of such stoichiometrical isomers could be observed populating up to $N_I \leq N_p$ sites with at least one proton? Mathematically speaking, this is a discrete number theory problem: This number is the integer partition of N_p , i.e. the number of ways to write N_p as sum of positive integers. For example, CH_4 , NH_3 , H_2O , HF , Ne represent only 5 out of all the 42 possible stoichiometrical isomers for $N_p = N_e = 10$. The total number of possible partitions corresponds to the partition function, which increases exponentially with N_p . The exponential increase, and an illustration of the emerging stoichiometries according to Young-Ferrers diagrams, are shown in Fig. 1. These degrees of freedom are rarely explored in nature except when it comes to radioactive decay, nuclear fusion, or nuclear synthesis in the early stages of our universe. Through interpolation, however, we can meaningfully render this space continuous, as illustrated using density functional theory (DFT) for the potential energies displayed in the inset of Fig. 1.

- (ii) *constitutional isomers*: “high” temperature regime. At high temperatures, only strong chemical bonding (covalent, ionic, metallic) survives. Corresponding Lewis structures enumerate many (but not all) of the possible constitutional isomers distinguishable as possible topologies, or molecular graphs, that can be constructed. The enumeration (and canonization) of all possible constitutional isomers has been the focus of long standing graph-theoretical efforts^{8–11}. The exponential scaling of their number is also evident for the recently published exhaustive list of small organic molecules¹². This is the regime in which isomerism occurs through conventional “chemistry”, i.e. reactions that lead from one constitutional isomer to another, usually under the influence of pressure, temperature, light, or in the presence of some catalytic agent. We can model this and the subsequent regimes using *ab initio* molecular dynamics methods⁴. Universal or reactive force-fields attempt to accomplish similar sampling^{13,14}.
- (iii) *conformational isomers*: “ambient” temperature. Folding and un-folding events, sampling of intramolecular degrees of freedom, for example around dihedral angles and similar processes take place at “ambient” temperatures. These isomers are typically sampled using force-fields that assume fixed molecular topologies and parameterized charges, dihedral and angular terms, in addition to the typical potentials used for the chemical bond, such as harmonic, Buckingham’s or Morse potentials.
- (iv) *weakly interacting systems*: “low” temperature. Supra-molecular assemblies, soft aggregates con-

dense to molecular liquids or solids. Typically modeled using classical effective Lennard-Jones type potentials.

In the remainder of this review we will discuss recent contributions that are consistent with *all* of the four regimes, accounting for all the spatial *and* elemental degrees of freedom $\{\mathbf{R}_I\}$ and $\{Z_I\}$. As also discussed below, we will ignore the electronic number as an independent variable since $N_e \approx N_p = \sum_I Z_I$ for most if not all possible scenarios.

III. MOLECULAR GRAND-CANONICAL ENSEMBLE

A. Theory

Much of conceptual density functional theory (DFT) concerns the energy response to infinitesimal variations in number of electrons N_e , or external potential, $v(\mathbf{r})$ ^{15,16}. While very important for interpreting orbitals, deriving reactivity indices, and even for redox-processes, the diversity (and combinatorial scaling) of CCS is rather due to variations in nuclear charge distribution, than due to variations in N_e . Consequently, for the following we will mostly be concerned with changes in nuclear charges. In order to offer a rigorous framework for explicit changes in $\{Z_I\}$, molecular grand-canonical ensemble DFT was introduced¹⁷, relying to a significant degree on preceding work^{18,19}. Only a brief summary is given here, for more details the reader is referred to the original contributions.

Assuming a classical nuclear charge distribution, $n_p(\mathbf{r})$, we can introduce an auxiliary grand-canonical variational energy functional for the aforementioned fictitious “very high” temperature regime (i),

$$\Omega[N_e, n_p(\mathbf{r})] = E[N_e, n_p(\mathbf{r})] - \mu_e \left(\int d\mathbf{r} n(\mathbf{r}) - N_e \right) - \mu_p \left(\int d\mathbf{r} n_p(\mathbf{r}) - N_p \right). \quad (1)$$

Where E, n, μ_e, μ_p correspond to the usual total potential energy functional, the electron charge density, and global electronic and nuclear chemical potentials, respectively. For high temperatures, entropy will prevail and the system would dissociate into H_{N_p} . For lower temperatures, the potential energy will dominate the free energy, and the energy of a single atom ($E(Z) \approx Z^{2.4} = N_p^{2.4}$) will dominate over the energy of many individual atoms that sum up to the same number of protons, $\sum_I Z_I^{2.4} \leq N_p^{2.4}$. Hence, the nuclear charge distribution would collapse onto a single site. For this discussion, we assume that the classical and fictitious self-repulsion of protons occupying the same nuclear site is switched off.

For the lower temperature regimes (ii)-(iv), the follow-

ing energy functional is more meaningful,

$$\Omega[N_e, n_p(\mathbf{r})] = E[N_e, n_p(\mathbf{r})] - \mu_e \left(\int d\mathbf{r} n(\mathbf{r}) - N_e \right) - \int d\mathbf{r} \mu_p(\mathbf{r}) \left(n_p(\mathbf{r}) - \sum_I Z_I \delta(|\mathbf{r} - \mathbf{R}_I|) \right) \quad (2)$$

where $\sum_I Z_I \delta(|\mathbf{R}_I - \mathbf{r}|)$ corresponds to the spatially resolved nuclear charge distribution. The nuclear chemical potential μ_p is no longer a global parameter but rather a locally defined Lagrange multiplier. Using an external potential that excludes the aforementioned intra-nuclear self-repulsion of protons (here through use of an error function), we find for the corresponding Euler equation,

$$\mu_p(\mathbf{r}) = \frac{\delta E}{\delta n_p(\mathbf{r})} = \sum_I \frac{Z_I \text{erf}[\sigma|\mathbf{R}_I - \mathbf{r}|]}{|\mathbf{r} - \mathbf{R}_I|} - \int d\mathbf{r}' \frac{n(\mathbf{r}')}{|\mathbf{r} - \mathbf{r}'|}, \quad (3)$$

—the electrostatic potential of the system. As such, starting with Q —a Legendre transformed energy functional of intensive properties μ_e and $\mu_p(\mathbf{r})$ —one can derive the Gibbs-Duhem equation for electrons *and* protons and electrons,

$$dQ[\mu_e, \mu_p(\mathbf{r})] = -N_e d\mu_e - \int d\mathbf{r} n_p(\mathbf{r}) \delta\mu_n(\mathbf{r}), \quad (4)$$

and obtain relationships between electronic hardness $\partial^2 E / \partial N_e$, molecular Fukui function $\partial \mu_p(\mathbf{r}) / \partial N_e$ ²⁰, and nuclear hardness, $\delta \mu_p(\mathbf{r}) / \delta n_p(\mathbf{r})$ ¹⁷. While the nuclear chemical potential is defined everywhere, its value at an atomic position quantifies the system’s first order energy response to a fractional change of the atom’s nuclear charge. Consequently, we dub $\mu_p(\mathbf{R}_I)$ the “alchemical potential” of atom I ¹⁹.

B. Interpolating pseudopotentials

Apart from radioactive processes alchemical changes obviously do not occur in reality. They offer, however, a rigorous mathematical way to render CCS continuous. Alchemical changes and potentials involving fractional nuclear charges are commonly used for two, often related, purposes: Either for the evaluation of free energy differences between different compounds, e.g. using thermodynamic integration²¹, $\Delta F = \int d\lambda \langle \partial E / \partial \lambda \rangle$; or for obtaining a set of gradients with dimension of N_I indicating the response of the system to a variation in nuclear charge on every site^{19,22}. In practice, we can calculate such changes through interpolation of nuclear charges in any basis set that is converged for all values of an interpolating order parameter, λ . For plane-wave pseudopotential implementations, the same can be accomplished by interpolation of pseudopotentials that replace the explicit treatment of the core electrons^{23–28}. The use of a plane-wave basis set is advantageous since it is independent of atomic

position and type, and will not introduce Pulay forces²⁹. The manipulation of pseudopotentials for affecting electronic structure properties is nothing out of the ordinary. It has successfully been deployed for an array of properties including relativistic effects³⁰, self-interaction corrections,^{31,32} exact-exchange and QM/MM boundary effects^{33,34}, van-der-Waals interactions^{35,36}, and widening the band gap^{37,38}. For fractional nuclear charges we can interpolate pseudopotentials, and evaluate properties as a function of order parameter, $0 \leq \lambda \leq 1$. An interpolation of pseudopotential parameters as a function of nuclear charge is shown in Fig. 2. Calculated properties as a function of such alchemical changes are illustrated in Figs. 1 and 3 for total potential energies, and protonation energies and polarizabilities, respectively. Note that the former property is not physical because of the arbitrary energy offset of pseudopotentials. This, however, is inconsequential, since most of chemistry deals with energy differences, and differences thereof. As shown in Fig. 3 for $\text{HCl} \rightarrow \text{NH}_3$, the use of pseudopotentials for alchemical changes can be particularly advantageous when it comes to transmuting elements from different rows of the periodic table while keeping constant the total number of *valence* electrons.

C. Free energy applications

Fractional charges were used to calculate free energy of solvation of ions in water³⁹. Sulpizi and Sprik rigorously explored the need for fractional nuclear charge calculations to obtain pK_a 's of various organic and inorganic acids and bases⁴⁰. In the case of free energy differences, fractional charges can also be avoided all together within a simple alternative and elegant interpolation scheme put forth by Alfè, Gillan and Price: Atomic forces are evaluated at both end-points ($\lambda \in 0, 1$), and λ dependent molecular dynamics trajectories are generated for atoms being propagated according to a linear combination of these forces using instantaneous λ values as weights⁴¹. This is to be compared to a trajectory that uses Hellmann-Feynman forces directly evaluated on the interpolated alchemical species, such as for the 0 Temperature limit of relaxing the geometry of a reaction barrier⁴². The limitations of Alfè's procedure are that (a) one requires twice as many self-consistent field calculations, namely for both end-points instead of a single one when using an alchemical interpolation (assuming of course that for both approaches the free energy integrand $\langle \partial E / \partial \lambda \rangle$, varies similarly with λ); and (b) that the number of atoms must be kept constant during the interpolation, significantly restricting the possible number of stoichiometries that can be explored. Both of these disadvantages can be avoided within the compound pair scheme discussed in section IV.

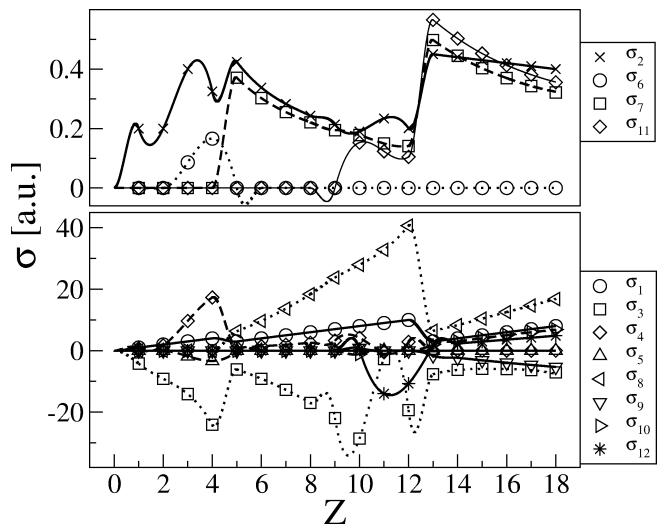


FIG. 2: Interpolation of Goedecker-Hutter pseudopotential parameters for the BLYP based DFT calculations^{43,44}. Parameters, shown as a function of nuclear charge, become polynomial regressions of third degree in Z for intervals of Z connected such that the sum is continuous and differentiable everywhere.

D. Design applications

Through use of a Taylor expansion, truncated after first order, we can also exploit the vector of alchemical potentials for gaining control over properties in the N_I -dimensional space of all the nuclei that are adjacent in the periodic table. Weigend, Schrodtr and Ahlrichs were probably the first to present such an application for the prediction of stability in binary atom clusters⁴⁵. Subsequently it was applied to drug binding energies within a QM/MM study¹⁹, demonstrated for hydrogen-bonded complexes, inter-converting methane, ammonia, water and hydrogen-fluoride while bound to formic acid²², and applied to the molecular Fukui function for tuning HOMO eigenvalues of boron/nitride derivatives of benzene²⁰. In Ref.⁴² this notion is exploited for the prediction of reaction barriers as well as oxygen adsorption energies on Pd_{79} derived core-shell metal nano-clusters that are catalyst candidates for the oxygen reduction reaction. The molecular Fukui function, in particular when evaluated at the position of the atom, was also discussed more recently by Cardenas et al.⁴⁶. Clearly, truncation of the Taylor expansion at second or higher order would be desirable to increase the accuracy of the alchemical predictions of the effect of atomic transmutations. Alas, higher order derivatives of the energy with respect to nuclear charges lead to computational overhead, they require the calculation of the perturbed electronic structure. Nevertheless, based on coupled perturbed self-consistent field theory the improved accuracy of higher order predictions was demonstrated very recently⁴⁷.

E. Other applications

Instead of varying the proton distribution $n_p(\mathbf{r}) = \sum_I Z_I \delta(|\mathbf{r} - \mathbf{R}_I|)$ in a compound, we can interpolate the external potential $v(\mathbf{r}) = \sum_I Z_I / |\mathbf{r} - \mathbf{R}_I|$ just as well. Albeit mixing up spatial and compositional degrees of freedom, this is more in line with conventional thought in quantum chemistry and conceptual DFT¹⁵ where the nuclear charge distribution is hardly mentioned explicitly. The route via the external potential has been pursued within the *Ansatz* of linear combination of atomic potentials in the research groups of Yang and Beratan⁴⁸, assigning atom-type specific weights to every atom site. Using simplified Hamiltonians, impressive results were obtained for the control of molecular hyperpolarizabilities^{49–53}, based on long-standing molecular design efforts for electronic properties well ahead of their time^{54,55}. This approach has also been adapted and explored for the purpose of crystal structure design using DFT⁵⁶. The functional second order derivatives with respect to external potentials have been published in Ref.⁵⁷. Analytical expressions for second order derivatives and linear response functions have very recently been proposed by Yang, Cohen, De Proft and Geerlings⁵⁸. The same authors also derived important constraints for the electronic structure that must be met by the exact exchange-correlation functional. In analogy to using constraints obtained for variable N_e , such as piecewise linear behavior and derivative discontinuities, in order to design improved density functionals^{59–61}, Cohen's current efforts are dedicated to variations in the external potentials that include fractional nuclear charges. The electronic structure for systems with $Z \rightarrow \infty$ has also been explored by Constantin et al.⁶².

IV. COMPOUND PAIRS

A. Background

The above discussed *Ansatz*, variational in a fractional nuclear charge distribution, defines an appealing, fully spatially resolved, index, i.e. a way to probe the sensitivity of a compound not only towards changes in any of its composing atoms but also with respect to adding new protons. However, for two reasons this approach can also be limited. First, severe constraints and preconceived insights are required to explore the N_I -dimensional space of all $\{Z_I\}$. Either because if Z_I is continuous it requires a bias potential towards integer numbers, possibly using a fictitious temperature, i.e. in analogy to the Fermi function for electrons. Or if Z_I is a combination of various atom-types, i.e. in line with the aforementioned linear combination of atomic potentials approach⁴⁸, the weight of one nuclear charge has to dominate so that it can safely be increased to 1, while decreasing all others to zero. Furthermore, constraints due to overall charge conservation, and electronic structure, have to be considered. For ex-

ample, consider an alchemical transmutation of $\text{H}_2\text{N-OH}$ into its iso-electronic stoichiometrical isomer hydrogen peroxide, HO-OH , through simultaneously and continuously decreasing and increasing by one the nuclear charge of the hydrogen and nitrogen atom, respectively. At some point of this conversion, the spin of the ground state surface will turn into a triplet surface, therefore requiring the consideration of *both* spin states along the interpolation path. Second, and more importantly, in order to carry out alchemical changes along columns in the periodic table, for example, a path following Z would have to fill up the shell to go through the entire period before one arrives at the desired elements. This implies significant variations in electronic configurations just to arrive at a target compound with a configuration likely to be very similar to the starting compound. For example, consider a system of 8 valence electrons, and Ne and Ar as starting and target compounds, respectively. Then, an iso-electronic path progressing with Z of the central atom, and saturating with hydrogens accordingly, would have to proceed through the following series of compounds, NaH_7 , MgH_6 , AlH_5 , SiH_4 , PH_3 , H_2S , and HCl , some of which not even likely to be covalently bound. Hence, while Taylor expansions in Z are quite predictive for adjacent elements—as mentioned in the preceding section—it is not surprising that their predictive power decays dramatically when it comes to predictions for changes up and down the columns in the periodic table. Obviously, matters only become worse when d - or f -elements have to be included, or when trying to make predictions by 2 or more rows down or upward.

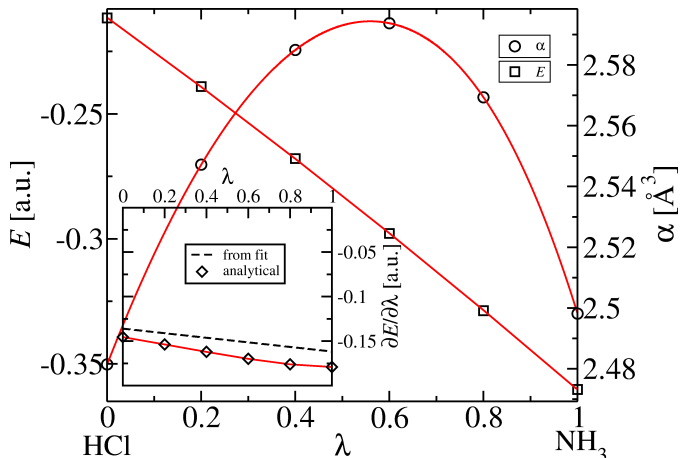


FIG. 3: Protonation energy, $E = E^{\text{X-H}^+} - E^{\text{X}}$, and static polarizability α of neutral species, as a function of order parameter λ driving $\text{X} = \text{HCl}$ into $\text{X} = \text{NH}_3$. The inset shows the derivative of the protonation energy, once evaluated analytically according to Hellman-Feynman in Eq. (8), and once from a quadratic fit to protonation energy. Heavy atoms of two end-point molecules were superimposed, their effective nuclear charge being $(1 - \lambda)7 + \lambda 5$, and hydrogens were placed in xy -plane. The protonating proton was placed in z , 1 Å above the heavy atom. All energies calculated with PBE functional and interpolated analytical pseudopotentials³⁰.

B. Theory

Albeit intuitive, the use of nuclear charges as interpolating variable is fortunately not mandatory. Instead, we can also use a generalized, and entirely arbitrary, interpolation procedure between any two pairs of compounds—as long as it is reversible and integrable, *any* path can be used to monitor any property that is a state function⁴. In all of the following, we will only consider interpolations between iso-electronic compounds, i.e. compounds with the same N_e in their Hamiltonian. As mentioned before, this is only a minor restriction because the diversity of CCS is by far not as much due to differences in N_e as it is due to differences in nuclear charge distribution. We can directly interpolate the nuclear charge distributions/external potentials/Hamiltonians of any two such iso-electronic compounds, A and B , e.g.

$$H(\lambda, \mathbf{r}) = H_A(\mathbf{r}) + \lambda(H_B(\mathbf{r}) - H_A(\mathbf{r})), \quad (5)$$

interpolating linearly and globally in order parameter λ . H_A and H_B denote the initial and final electronic Hamiltonian of the two compounds, defining the corresponding boundary conditions $H(\lambda = 0)$ and $H(\lambda = 1)$, respectively. For any iso-electronic Hamiltonian linear in λ , the potential energy is not necessarily linear. In fact, the electronic potential energy of a linearly interpolated Hamiltonian will always be concave, i.e. equals or larger than a straight line connecting the energies of compound A and B , $E^{lin} = E_A - \lambda(E_B - E_A)$. This inequality follows from the variational principle and can easily be shown: Eq. 5 implies,

$$E[H(\lambda)] = \langle H(\lambda) \rangle_\lambda = E_A[n_\lambda] + \lambda(E_B[n_\lambda] - E_A[n_\lambda]), \quad (6)$$

where $\langle \dots \rangle_\lambda$ now correspond to the usual quantum mechanical Bra-Ket notation, denoting the expectation value with the wavefunction, or the density functional (in an orbital-free exact DFT world), evaluated for the Hamiltonian at λ , i.e. $E[H(\lambda)] = \langle \Psi_\lambda | H(\lambda) | \Psi_\lambda \rangle = \langle H(\lambda) \rangle_\lambda = E_\lambda[n_\lambda]$. $E_A[n_\lambda]$ and $E_B[n_\lambda]$ denote the energies of compound A or B evaluated using the wavefunctions (or density in the case of orbital free DFT) obtained at λ . Note that $n_{\lambda=0} = n_A$, and $n_{\lambda=1} = n_B$. Subtracting $E^{lin}(\lambda)$ and regrouping yields,

$$\begin{aligned} E[H(\lambda)] - E^{lin}(\lambda) &= (E_A[n_\lambda] + \lambda(E_B[n_\lambda] - E_A[n_\lambda]) \\ &\quad - (E_A[n_A] + \lambda(E_B[n_B] - E_A[n_A]))) \\ &= (1 - \lambda)(E_A[n_\lambda] - E_A[n_A]) \\ &\quad + \lambda(E_B[n_\lambda] - E_B[n_B]), \\ &\geq 0. \end{aligned} \quad (7)$$

where the prefactors of the energy differences $\lambda, (1 - \lambda) \geq 0$ by definition, and where $E_A[n_A] \leq E_A[n_\lambda]$ and $E_B[n_B] \leq E_B[n_\lambda]$ because of the variational principle. Consequently, analogous inequalities will hold for any property for which there is a variational principle,

e.g. also for the polarizability due to Pearson's maximum hardness principle⁶³. This inequality is on display for the static polarizability, fractionally transmutating a hydrogen chloride molecule into ammonia (Fig. 3). Similarly, potential energy inequalities in between different molecules were proposed by Mezey in the eighties⁶⁴.

Analytical first order derivatives of the energy as a function of some iso-electronic change in the Hamiltonian can easily be calculated using the Hellmann-Feynman (HF) theorem⁶⁵, as proposed and demonstrated for HOMO eigenvalues in Ref.⁶⁶, $\partial E / \partial \lambda = \langle \partial H / \partial \lambda \rangle_\lambda$. For a linearly interpolating Hamiltonian, such as in Eq. (5), this leads to,

$$\begin{aligned} \frac{\partial E[H(\lambda)]}{\partial \lambda} &= \langle H_B - H_A \rangle_\lambda = \int d\mathbf{r} n_\lambda(\mathbf{r}) \times (v_B(\mathbf{r}) - v_A(\mathbf{r})) \\ &= \langle H_B \rangle_\lambda - \langle H_A \rangle_\lambda = E_B[n_\lambda] - E_A[n_\lambda]. \end{aligned} \quad (8)$$

The protonation energy, and its derivative, also feature in Fig. 3 for the same transmutational change, $\text{HCl} \rightarrow \text{NH}_3$. As mentioned before, the use of pseudopotentials/valence electron densities fortunately renders straightforward the application of the HF theorem according to Eq. (8) even for changes that involve elements from differing rows in the periodic table.

Thermodynamic integration of $\partial E / \partial \lambda$ over λ yields any properties related to free energy differences. In the case of compound design the approach is slightly different, we would like to expand the energy of a new compound B in terms of a reference compound A and its derivatives,

$$E_B \approx E_A[n_A] + \frac{\partial E_A[n_A]}{\partial \lambda} \Delta\lambda + \frac{1}{2} \frac{\partial^2 E_A[n_A]}{\partial \lambda^2} \Delta\lambda^2 + \text{HOT}, \quad (9)$$

HOT standing for higher order terms, and $\Delta\lambda = 1$. Unfortunately, when making predictions with a linearly interpolated Hamiltonian, the first order derivative term according to Eq. (8) is not necessarily predictive⁶⁶. Unfortunately, the inclusion of higher order derivatives in Eq. (9) might not only improve the prediction, as found for statistical mechanical averages⁶⁷, but it also requires the evaluation of the perturbed wavefunction, e.g. through the use of linear response theory^{33,68}, thereby defying the original purpose of predicting a new compound's energy without having to calculate its wave function. Nevertheless, for external potentials this has been carried out within conceptual DFT⁵⁷, and very recently even analytically⁵⁸.

In order to improve the predictive power of the first order term in Eq. (9), an empirical correction has been introduced that "linearizes" the energy through a global yet non-linear Hamiltonian, $H(\lambda) = H_A + f_{AB}(\lambda)(H_B - H_A)$ ⁶⁶. If we assume $f(\lambda)$ to be a second order polynomial in λ , two coefficients are determined by the boundary conditions that $H(\lambda = 0) = 0$, and $H(\lambda = 1) = 1$, leaving one additional degree of freedom. We can obtain the remaining degree of freedom as a parameter from an

arbitrary second iso-electronic compound pair, CD , such that the energy becomes linear in λ . The resulting expansion up to first order in Eq. (9) then becomes,

$$E_A \approx E_A[n_A] + C_{CD}^{ref} \frac{\partial E_A[n_A]}{\partial \lambda} \Delta \lambda. \quad (10)$$

Here, C_{CD}^{ref} is the ratio between the energy difference and the HF derivative of the additional reference compound pair, C and D (its Hamiltonian being linearly interpolated), and $\partial_\lambda E_A[n_A]$ is determined according to Eq. (8) as $E_B[n_A] - E_A[n_A]$. This bears resemblance to a long tradition in physical chemistry, namely the use of reference compounds for electrode potentials or enthalpies.

The idea to use alternative, non-linear, interpolations is not new within the molecular mechanics research. In the context of electronic structure theory non-linear alchemical paths were also explored for chemical binding⁶⁹, and nuclear quantum effects⁷⁰. Various open questions deserve further investigation, such as transferability and choice of reference coefficients, iso-electronic changes using valence electrons only versus all electron description, non-iso-electronic changes, necessary accuracy when providing the input of target compound B , i.e. also its geometry, ionic forces of B , etc. These answers are likely to depend on systems and properties.

C. Control of ligand binding

In this section, we exemplify the use of the reference coefficients [Eq. (10)] for increasing the predictive power of the HF derivatives of linearly interpolated Hamiltonians. We refer to state-of-the-art van der Waals corrected DFT^{35,71} to accurately estimate interaction energies with binding targets across CCS. We will consider a small yet illustrative set of mutants of the ellipticine molecule. Ellipticine is a naturally occurring anti-cancer drug with various binding targets. As also illustrated in Fig. (4), its dominant mode of binding to DNA is intercalation. Structural data as well as studies on drug analogues are readily available^{72,73}. We will probe the versatility of the linearizing scheme for controlling ellipticine-derivatives/DNA binding⁷⁴. Clearly, for the eventual control of ligand binding the property of interest is not the potential energy of interaction but rather the free energies of binding: Solvation or entropic contributions can be crucial, as is well known in general⁷⁵, and in the particular case of ellipticine⁷⁶. For example, Tidor⁷⁷, and Oostenbrink and van Gunsteren^{78,79} have carried out similar work in the sense of interpolating ligand candidates, by calculating free energies of binding, and using molecular force-fields. However, here we will focus on the potential energy of interaction. Subsequent work in the future can deal with the inclusion of thermal and solvent effects for instance using *ab initio* molecular dynamics techniques⁸⁰ in conjunction with QM/MM⁸¹ calculations. Moreover, even at the mere potential energy

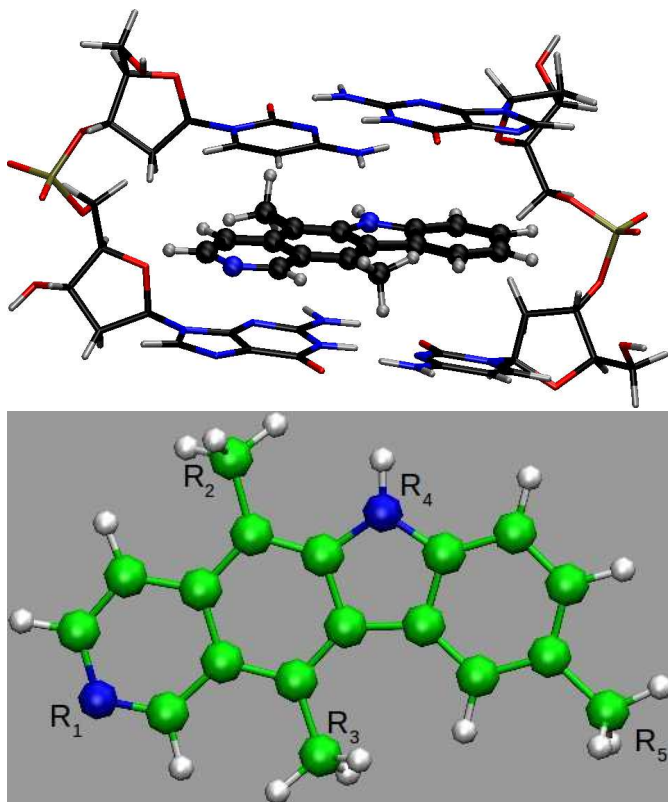


FIG. 4: TOP: Cluster model of drug intercalated in between two Watson-Crick base-pairs connected by sugar puckers and phosphate groups. BOTTOM: Neutral “wild-type” ellipticine, R_i denote sites of groups permitted to mutate (see TAB. I).

electronic structure level of theory, the accurate quantification and control of intercalated ellipticine derivatives is challenging: vdW forces dominate the binding. Recent studies have already explored the binding of ellipticine and how its vdW forces can be accounted for at the employed electronic structure level^{82–84}.

Let us consider the intercalation energy for the complex depicted in Fig. (4) for mutations at the five sites indicated in the bottom panel. In analogy to protein or DNA sequences, an (arbitrary) relevant subspace of CCS is defined in TAB. (I) as a matrix that corresponds to an alphabet of iso-electronic (in valence electron number) functional groups at each of the selected sites. Note that variation in molecular combinations of letters of this alphabet are capable to not only revert dipole-moments, they can also act either as hydrogen bond acceptors (lone pair in OH/Cl) or donors (NH₂, proton in OH). Clearly, the alphabet can easily be extended to accommodate further effects, for example with electron donating/withdrawing or hyperconjugating groups etc. Conformational degrees of freedom can be encoded explicitly, as it is done for the hydroxyl groups in TAB. (I).

Within this restricted CCS, any given molecule is represented by the sequence of functional groups distributed

TABLE I: Exemplary alphabet for mutants of ellipticine as oriented in Fig. 4, defining a CCS with $4 \times 6^4 = 5184$ molecules. Highlighted in red are all functional groups whose mutations have been considered. Predictions are displayed in Fig. 5. The “wild-type” ellipticine drug is encoded as (21121) with three functional groups coming from the first column, and the functional groups at site R_1 and R_4 coming from the second column.

site vs. group	1	2	3	4	5	6
R_1	CH	N	SiH	P	-	-
R_2	CH ₃	NH ₂	OH ^{left}	OH ^{right}	F	Cl
R_3	CH ₃	NH ₂	OH ^{left}	OH ^{right}	F	Cl
R_4	CH ₂	NH	O	SiH ₂	PH	S
R_5	CH ₃	NH ₂	OH ^{left}	OH ^{right}	F	Cl

over the five sites. For example, the “wild-type” ellipticine in Fig. (4) would correspond to (21121), *i.e.* 2 for N at R_1 , 1 for CH₃ at R_2 , R_3 , and R_5 , and 2 for NH₂ at the R_4 . Let us exemplify a DFT+vdW based prediction of the binding energy of another mutant: (21121) is predicted to bind to the DNA cluster in Fig. 4 with $E_{21121} = -38.5$ kcal/mol⁸². For predicting the single point mutation (21121)→(21125) (changing CH₃ into F at R_5), one would have to predict a target value of $E_{21125} = -36.9$ kcal/mol. The derivative based prediction according to first order term in Eq. (9) is calculated to be, $E_{21125} \approx E_{21121} + \partial_\lambda E_{21121} = -38.5 + 1.4 = -37.1$ kcal/mol. Inclusion of reference coefficient [Eq. (10)], and using compound pair (11121)/(11125) as a reference, yields $E_{21125} \approx E_{21121} + C_{ref} \times \partial_\lambda E_{21121} = -38.5 + 1.3 \times 1.4 = -36.7$ kcal/mol.

In order to gain a more representative idea of the predictive power of this method, Fig. 5 features the outcome for a small subspace of the CCS highlighted in red in TAB. I: Eight compounds have been considered involving permutations at R_1 , R_4 , and R_5 , each with two possible functional groups. Predictions based on all the possible derivatives among these compounds, with and without reference coefficients (as obtained from compound pairs not involved in the transmutation), are compared to calculated binding energies. Despite the several outliers that deviate substantially, the use of reference compounds dramatically improves the overall prediction.

For comparison, we also include predictions based on the additive assumptions that the influence of the rest of the molecule cancels when considering an interpolation for the same pair of functional groups. Specifically, we estimate the binding energy of B simply by adding the difference in binding energy of a reference compound pair, CD , to the binding energy of A,

$$E_B \approx E_A + \Delta E = E_A + (E_D - E_C) \quad (11)$$

As shown in Fig. 5, also this prediction yields remarkable good correlation—with less pronounced outliers. Analysis of the distribution of errors, however, suggests that in spite of the outliers the normal distribution of predictions around the ideal correlation is superior for predic-

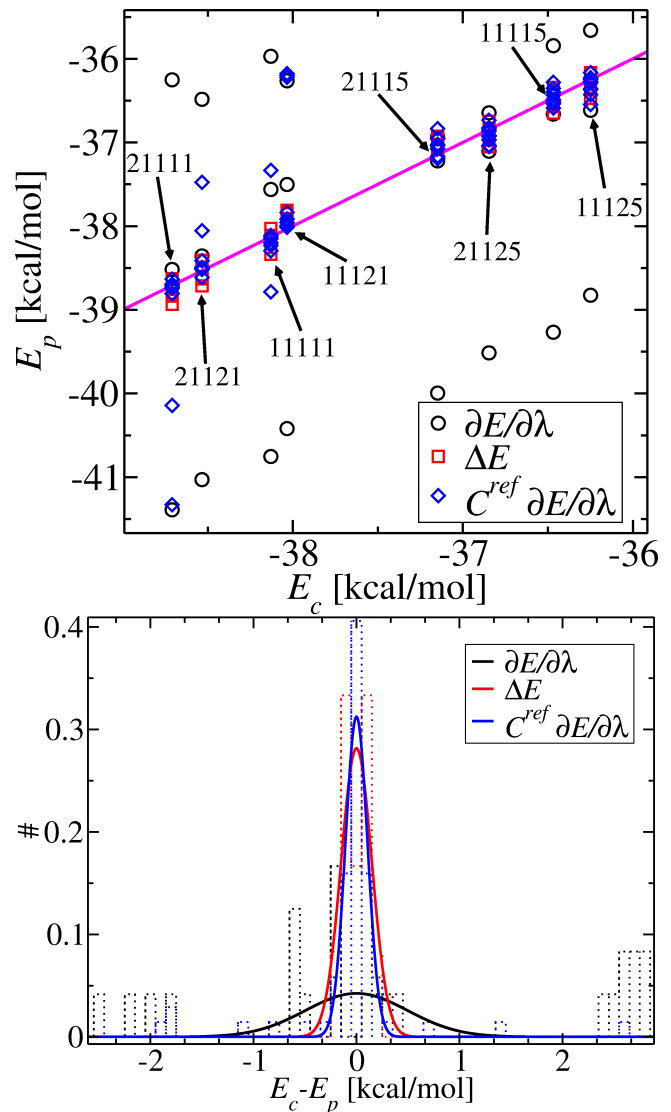


FIG. 5: TOP: Correlation for eight compounds from alphabet in TAB. I. Predictions made using first order derivatives only ($\partial E/\partial\lambda$, Eq. (9)), energy difference of reference compound pairs (ΔE , Eq. (11)), or $C^{ref} \partial E/\partial\lambda$, Eq. (10). BOT-TOM: Normalized histogram and corresponding normal distribution of error over 72 predictions, $\Delta_{c-p} E^{int} = \text{calculated } E^{int} - \text{predicted } E^{int}$.

tions made with the product of derivative and reference coefficient (Bottom of Fig. 5).

D. Win a prize

The numerical illustrations in the previous section, as well as in Ref.⁶⁶, suggest that efforts to linearize the property through use of alternative, non-linear interpolations of the Hamiltonians are worthwhile. Strictly speaking, however, due to the use of reference compound pairs, the aforementioned interpolation constitutes no longer a first

principles but rather an empirical and heuristic *Ansatz*. What is needed instead, is an *ab initio* interpolating procedure that linearizes the energy (or other properties) in order parameter, such that the first order Taylor expansion based on the HF derivative is sufficiently accurate to predict properties of other compounds².

Inspired by Erdős’ habit to offer cash awards for solutions to outstanding mathematical problems, the author has thus decided to offer the equivalent of an ounce of gold to the first person who presents an *ab initio* solution to this problem. Specifically, the challenge reads: Find—or show non-existence of—an *ab initio*, i.e. valid for *any* external potentials, interpolating transform $f_{AB}(\lambda)$ for which two different but iso-electronic molecular Hamiltonians with energies E_A and E_B interconvert such that the electronic ground state potential energy $E = \langle H(f(\lambda)) \rangle$, is linear in order parameter λ , and that consequently the HF derivative is given by,

$$\left. \frac{\partial E(\lambda)}{\partial \lambda} \right|_{\lambda} = \left\langle \frac{\partial H(f(\lambda))}{\partial \lambda} \right\rangle_{\lambda} = E_B - E_A. \quad (12)$$

Here, $0 \leq \lambda \leq 1$, and $E(\lambda = 0) = \langle H(f(\lambda = 0)) \rangle = \langle H_A \rangle = E_A$, and $E(\lambda = 1) = \langle H(f(\lambda = 1)) \rangle = \langle H_B \rangle = E_B$. Further details can be found in footnote⁸⁵.

We can exemplify the challenge by solving it for the non-relativistic hydrogen-like single atom with only one electron. In this case, $E(\lambda) = aZ(\lambda)^2$, where a is a constant, and the nuclear charge Z is a function of interpolating parameter λ ¹³⁹. For an interpolation linear in the Hamiltonian, $Z(\lambda) = Z_A + \lambda(Z_B - Z_A)$, and the energy is therefore clearly quadratic in λ . The desired behavior of a linearized energy would be,

$$\begin{aligned} E^{lin}(\lambda) &= E(Z_A) + \lambda(E(Z_B) - E(Z_A)) \\ &= a(Z_A^2 + \lambda(Z_B^2 - Z_A^2)). \end{aligned} \quad (13)$$

Equating this to $aZ(\lambda)^2$ and solving for $Z(\lambda)$ yields the corresponding interpolating function:

$$Z(\lambda) = \sqrt{Z_A^2 + \lambda(Z_B^2 - Z_A^2)}. \quad (14)$$

As suggested above in the challenge, we can test this interpolation to confirm if indeed we find the desired slope for the linearized energy, $E_B - E_A$. Application of the chain rule, and insertion and differentiation of Eq. (14) confirms the expected result,

$$\begin{aligned} \left. \frac{\partial E}{\partial \lambda} \right|_{\lambda} &= \left. \frac{\partial E}{\partial Z} \frac{\partial Z(\lambda)}{\partial \lambda} \right|_{\lambda} = 2aZ(\lambda) \frac{Z_B^2 - Z_A^2}{2\sqrt{Z_A^2 + \lambda(Z_B^2 - Z_A^2)}} \Big|_{\lambda} \\ &= a(Z_B^2 - Z_A^2) \equiv E_B - E_A. \end{aligned} \quad (15)$$

As such Eq. (14) linearizes the energy in λ . The challenge of the prize consists of finding an analogous expression for molecules, i.e. a spatially resolved and λ dependent distribution of nuclear interpolations, $\{Z(\mathbf{R}_I, \lambda)\}$, that

drive all atoms in compound A to atoms in compound B while linearizing the potential energy.

Note that a naive extension of Eq. (14) to assemblies of atoms,

$$\frac{\partial E}{\partial \lambda} = \sum_{I \in A, B} \frac{\partial E}{\partial Z_I} \frac{\partial Z(\lambda)}{\partial \lambda} = \sum_{I \in A, B} \mu_p(\mathbf{R}_I) \frac{\partial Z(\mathbf{R}_I, \lambda)}{\partial \lambda} \quad (16)$$

does not constitute a practical approximate solution to the challenge. $\mu_p(\mathbf{R}_I)$ denotes the “alchemical” potential mentioned above which corresponds to the electrostatic potential at \mathbf{R}_I without the repulsion due to Z_I . μ_p will not necessarily cancel the square root term in the denominator of the derivative in Eq. (14), which consequently diverges if λ and $Z_A(\mathbf{R}_I)$ equal 0.

V. STATISTICAL METHODS

A. Inductive reasoning from first principles

Within statistical mechanics the numerical prediction of macroscopic observables from atomistic simulation requires repeatedly calculating microscopic states, using electronic structure theory, atomistic or coarse-grained force fields, and averaging in an appropriate ensemble. Philosophically speaking, the exercise of performing such computational “experiments” is an application of *deductive* reasoning to increase knowledge. But also when exploring CCS in terms of ensembles of potential energy hypersurfaces by repeatedly solving SE for N different compounds deductive reasoning is at work. Since the size of CCS is prohibitively large, its exhaustive exploration through screening with SE is impossible. While some interpolating λ schemes use statistical mechanics for a preselected set of compounds^{77,78}, a rigorous way to systematically and generally gain quantitative insights is desirable. This task can be accomplished through the application of *inductive* reasoning.

Historically, the role of inductive reasoning for chemistry is considerable, Mendeleev’s table, the Hammett equation, or Pettifor’s structure maps^{86–88} are all based on inferred phenomenological relationships. Further examples include widely spread rules and notions of chemistry, such as the chemical bond, atomic charges, or aromaticity. While popular and useful to the experimental chemist conventional quantum chemistry, based on deductive reasoning, is still struggling to account for these notions⁸⁹. Recent advances in statistical data analysis, methods^{90–93} and applications in other areas of science and engineering, such as searching the internet, automated locomotion (self-driving cars), algorithmic trading, or brain-computer interfaces, strongly suggest that they will also play an increasingly important role in chemistry. Examples of first efforts to quantitatively infer laws for atomistic simulations include “Learning On The Fly”⁹⁴, or “force-matching”^{95,96}. More sophisticated

statistical learning methods have been applied to the training of exchange correlation functionals in density-functional theory^{97,98}, or to parameterizing interatomic force fields^{99–104}. Support vector machines have been shown to quantify basis-set incompleteness¹⁰⁵. Gaussian kernel based machine learning (ML) for very accurate reactive force-fields was introduced by Bartok et al.¹⁰⁶. Contributions by Curtarolo, Hautier, and Ceder combine data-mining with mean-field electronic structure theory^{107–109}. Even the learning of reorganization energies that enter Marcus charge transfer rates are promising^{110,111}. Very recently, kernel based ML models have also delivered promising results for learning electron density functionals¹¹², or transition state theory dividing surfaces that determine reaction rates¹¹³. Bayesian error estimates and cross-validation methods have also been applied to the development of exchange-correlation models with controlled transferability¹¹⁴.

Within the bioinformatics and cheminformatics communities the development of quantitative structure property relationships (QSPRs) has a long tradition. QSPRs, relying on similar statistical frameworks (ML, cross-validated training, principal component analysis, etc.), deliberately attempt to circumvent solving the underlying laws of physics by directly correlating system parameters (descriptors) with macroscopic properties of interest. Conventionally, QSPRs are based on descriptors that explicitly forsake atomic resolution in the first principles sense. A large variety of such QSPR descriptors for various properties has been proposed^{115–118}. Two such descriptors, the molecular signature by Faulon and coworkers¹¹⁹ and a combination of HOMO eigenvalues of charged and neutral species, have recently yielded promising results for the QSPR modeling of a first principles property, the reorganization energy, in the CCS of polycyclic aromatic hydrocarbons (PAHs)¹¹¹. PAHs are used in discotic liquid crystals which self-assemble into columnar liquid crystal structures, implying their usefulness for organic photo-voltaic applications¹²⁰.

In this section we will discuss the application of *ab initio* statistical learning approaches to previously obtained first principles data for N compounds. Merely based on the data, QSPRs can be inferred, i.e. “learned”, and subsequently be used to avoid the cumbersome task of having to explicitly model all the underlying physical degrees of freedom of electrons and nuclei. As such, ML estimates solutions of SE for a new, i.e. “unseen”, molecule B simply by evaluating an analytical expression $E^{est}(B)$ that (explicitly or implicitly) encodes the data of N other molecules. Obviously, any such inferred relationships are inherently limited in accuracy by the quality of the data used for training.

B. Machine learning in CCS: The quantum machine

Recently, a kernel ridge regression approach to learn DFT atomization energies across CCS has been introduced¹²¹. Unlike ordinary QSPR approaches, this ML model is free of any heuristics. It exactly encodes the supervised learning problem posed by SE, i.e. instead of finding the wavefunction Ψ which maps the system’s Hamiltonian to its energy, $H(\{Z_I, \mathbf{R}_I\}) \xrightarrow{\Psi} E$, it directly maps system to energy based on N examples given. In the limit of converged N , i.e. sufficiently dense system coverage, the ML model is therefore a formally exact inductive equivalent to the deductive solution of SE through use of approximate wave-functions (such as separability of nuclear and electronic wavefunction or single slater determinants), Hamiltonians (such as certain exchange-correlation potentials), and self-consistent field procedure to minimize the energy. In Ref.¹²¹ numerical evidence is given for this idea. Specifically, for a diverse set of organic molecules, one can show that a ML model can be used instead, $\{Z_I, \mathbf{R}_I\} \xrightarrow{\text{ML}} E$. After training, solutions to SE can be inferred for out-of-sample, i.e. “unseen”, compounds that differ either in geometry or in composition or in both. The evaluation of an estimate is ordinarily negligible in terms of computational cost, i.e. milli seconds instead of hours on a conventional CPU, while yielding an accuracy competitive with the deductive approaches of modern electronic structure theory. As within any inductive approach, the accuracy is limited by the domain of applicability as defined by the data used for training, i.e. robust results can only be expected in interpolating regimes with sufficient coverage. Within the Gaussian kernel model, the energy of a query molecule \mathbf{M}_A is given as a sum over N molecules in the training set,

$$E^{est}(\mathbf{M}_A) = \sum_{i=1}^N \alpha_i e^{-\frac{d(\mathbf{M}_A, \mathbf{M}_i)^2}{2\sigma^2}}. \quad (17)$$

Each training molecule i contributes to the energy according to its specific weight α_i , scaled by a Gaussian in its distance to \mathbf{M}_A , $d(\mathbf{M}_A, \mathbf{M}_i)$. For given length-scale σ and regularization parameter λ , $\{\alpha_i\}$ are obtained by solving the regression problem,

$$\min_{\alpha} \sum_i (E^{est}(\mathbf{M}_i) - E_i^{ref})^2 + \lambda \sum_i \alpha_i^2. \quad (18)$$

σ and λ are hyperparameters. This regularized model limits the norm of regression coefficients, $\{\alpha_i\}$, thereby improving the transferability of the model to new compounds. All regression coefficients and hyper-parameters are determined by cross-validation on data stratified training sets^{92,93}.

So far this model has been trained and validated only in its most rudimentary form for atomization energies of a small set of interesting compounds. Specifically,

molecular atomization energies at the hybrid DFT level of theory^{5,122–125} have been used for training on up to $N \approx 7000$ molecules from the GDB data base¹² (see Fig. 6 for an illustration), for which mean absolute errors of less than 10 kcal/mol have been obtained. The choice of hybrid DFT is motivated by relatively small errors (< 5 kcal/mol) for thermo-chemistry data that includes molecular atomization energies¹²⁶. While 10 kcal/mol is still far from “chemical accuracy” (≈ 1 kcal/mol), more recent progress has not only led to atomization errors with less than 3 kcal/mol accuracy¹²⁷, but also includes other electronic properties, such as frontier eigenvalues, polarizability, and excitation energies¹²⁸.

An appealing advantage of analytical models, independent if obtained from physical insight or statistical regression, is their amenability to physical analysis. For example, unlike electronic structure methods, otherwise ill-defined concepts such as distance/neighborhood/similarity in CCS can now be quantified within the “world” of the ML model. Specifically, Eq. (17) gives the energy of a query molecule \mathbf{M}_A as an expansion in compound space spanned by reference molecules $\{\mathbf{M}_i\}$: The regression weights $\{\alpha_i\}$ are scaled by the similarity between query and reference compound as measured by a Gaussian of the distance. Hence, α_i assigns a positive or negative weight to molecule i . Within the compound space used as reference, molecules therefore can be ranked according to their $|\alpha|$. However, since $\{\alpha_i\}$ are regression coefficients in a non-linear model, i.e. after a non-linear transformation of the training data, the resulting energy contributions are specific to the employed training set without general implications for other properties or regions of compound space. The locality of the model is measured by σ , enabling the definition of a critical distance of locality, d_c , i.e. only if $d(\mathbf{M}_A, \mathbf{M}_i) \leq d_c$ will \mathbf{M}_i contribute to the energy of \mathbf{M}_A more than some threshold energy E_c . Rearranging summands in Eq. (17) leads to $d_c(\mathbf{M}_A, \mathbf{M}_i) = \sigma \sqrt{2 \ln[\alpha_i/E_c]}$. For atomization energies, and the chemical space considered in Ref.¹²¹, i.e. with a critical distance ≤ 400 Bohr (see TOP of Fig. 6) the ML results suggest that the model becomes local when $\sigma \leq 60$ Bohr, for the average α , and for $E_c = 1$ kcal/mol. Such σ values are achieved when the number of molecules in training set N exceeds ~ 5000 . In other words, for $N \leq 5000$, the model is global, i.e. all reference compounds contribute with more than 1 kcal/mol to any prediction made. See BOTTOM of Fig. 6 for the N dependence of σ and λ .

C. Coulomb matrix descriptor

To represent compounds, a wide variety of “descriptors” is in use by statistical methods for chem- and bio-informatics applications^{115–119}. The descriptor introduced by Rupp et al.¹²¹ is based solely on coordinates and nuclear charges, and dubbed “Coulomb-matrix”, \mathbf{M} ,

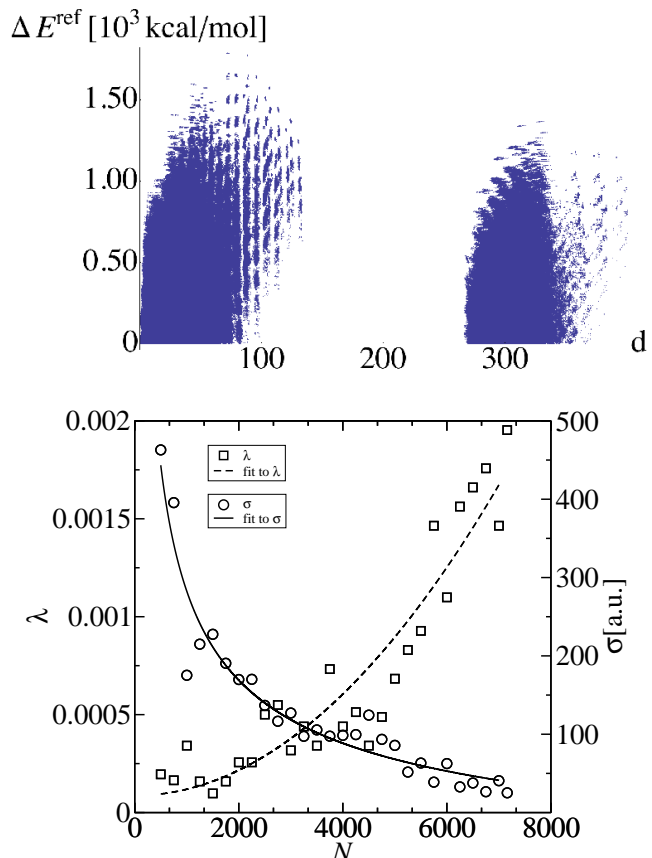


FIG. 6: TOP: Distance distribution in GDB-13 for the 7000 smallest molecules, $\Delta E^{\text{ref}} = |E_i - E_j|$ versus $d(\mathbf{M}_i, \mathbf{M}_j)$. BOTTOM: N dependence of σ and λ .

a symmetric square matrix of $N_I \times N_I$ dimensions,

$$M_{IJ} = \begin{cases} 0.5 Z_I^{2.4} & \forall I = J, \\ \frac{Z_I Z_J}{|\mathbf{R}_I - \mathbf{R}_J|} & \forall I \neq J. \end{cases} \quad (19)$$

The diagonal elements, $E(Z_I) \approx 0.5 Z_I^{2.4}$, correspond to a polynomial fit to free atom energies¹²⁹. The off-diagonal elements correspond to the Coulomb repulsion between atoms I and J . For a data set containing molecules with differing number of atoms, all the $\{\mathbf{M}\}$ of all the smaller systems are extended by zeros until they reach the dimensionality of the largest molecule in the training set. The Coulomb-matrix can easily be extended to account for extended or condensed phase systems: Let N_J be the number of atoms in the unit cell, and let N_I be the number of atoms in unit cell plus sufficiently large surrounding environment, then define M_{IJ} as above except that all off-diagonal elements are set to zero for all I and J larger than N_J .

We can measure the distance between two molecules by the Euclidean norm of their diagonalized Coulomb matrices: $d(\mathbf{M}_A, \mathbf{M}_B) = d(\epsilon_A, \epsilon_B) = \sqrt{\sum_I |\epsilon_{I \in A} - \epsilon_{I \in B}|^2}$, where ϵ are the eigenvalues of \mathbf{M} in order of decreasing absolute value. The physical meaning of represent-

ing CCS in this way can easily be understood by considering the simplest of all molecules, homo-nuclear diatomics (i.e. $Z = Z_1 = Z_2$ and $r = |\mathbf{R}_1 - \mathbf{R}_2|$). Any corresponding \mathbf{M} is then simply defined by its two eigenvalues, the roots of its characteristic polynomial, $\epsilon_{1/2} = 0.5Z^{2.4} \pm Z^2/r$. When measuring similarity between two such diatomics with different interatomic distances, r_A and r_B , the measure of similarity reduces to $d(\mathbf{M}_A, \mathbf{M}_B) = \sqrt{2} Z^2 (r_B - r_A)/(r_A r_B)$; and the corresponding estimated potential energy curve for any new interatomic distance, r_A , as trained on N other interatomic distances, $\{r_i\}$, is given by

$$E^{est}(r_A) = \sum_{i=1}^N \alpha_i \exp \left[-\frac{Z^4 (r_i - r_A)^2}{\sigma^2 r_A^2 r_i^2} \right]. \quad (20)$$

In complete analogy, a ML model of the homo-nuclear dimer can also analytically be understood in terms of other homo-nuclear dimers with differing atomic numbers, hetero-nuclear dimers, or hetero-nuclear trimers. The ease of differentiation with respect to not only geometry ($\partial_r E^{est}$) but also with respect to composition ($\partial_Z E^{est}$) illustrates further advantages of such a simple model.

The Coulomb matrix uniquely encodes any compound because stoichiometry as well as atomic configuration are explicitly accounted for. Even homometric molecules¹³⁰, see Fig. 7, are uniquely encoded by \mathbf{M} . Symmetrically equivalent atoms will contribute equally, and the representation is rotationally and translationally invariant. In order to gain invariance of \mathbf{M} with respect to the index ordering of atoms one can either diagonalize, sort rows and columns according to their norm, or use sets of matrices with permuted rows and columns. Using the eigenvalues of \mathbf{M} will yield an undercomplete representation. As with any coarsened representation, the N_I degrees of freedom represented by eigenvalues will fail to uniquely represent the full set of $3N_I - 6$ degrees of freedom for any non-linear molecule with more than three atoms^{131,132}. While sorting by the norm of rows (or columns) leads to an overcomplete, index invariant, and unique representation, the matrix is no longer differentiable for any combination of matrix entries that could be achieved through changes in geometry or in nuclear charges. Extending the representation by randomly permuted variants of Coulomb matrices is feasible, and leads to dramatic improvement in predictive accuracy^{127,128}. To encode known invariances through such data extension has also been successful for improving the accuracy of handwritten digit recognition¹³³. Due to disadvantageous scaling, this approach might prove problematic, however, when it comes to larger systems. As discussed in Ref.¹³⁴, these are all crucial criteria for representing atomistic systems within statistical models.

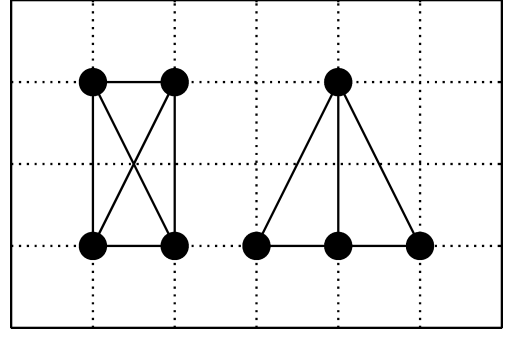


FIG. 7: Sketch of two homometric molecules (same stoichiometry, same sum of interatomic distances) from Ref.¹³⁵. The Coulomb-matrix (sorted or a set of its permutants) can distinguish these two molecules^{130–132}.

D. Alternative descriptors for CCS

We shall now discuss more sophisticated alternatives to the Coulomb-matrix. An intuitive extension is to assume a matrix with an interatomic potential form. This could be worth-while as long as the incurred computational overhead is small by comparison to the method used to generate the reference data. For example,

$$M_{IJ}^{LJ} = \begin{cases} 0 & \forall I = J, \\ \epsilon_{IJ} \left(\left(\frac{r_{IJ}^{eq}}{r_{IJ}} \right)^{12} - 2 \left(\frac{r_{IJ}^{eq}}{r_{IJ}} \right)^6 \right) & \forall I \neq J. \end{cases} \quad (21)$$

would correspond to the Lennard-Jones analog to the Coulomb-matrix. Similarly, a Morse or Buckingham matrix could be constructed. One could even conceive to go beyond such pair-wise approaches and introduce interatomic 3 and higher order terms in the form of molecular tensors. But also electronic structure models can be encoded in terms of such a representation, such as extended Hückel theory, semi-empirical quantum chemistry or tight-binding models. For example, an orbital free Thomas-Fermi DFT representation¹³⁶ is possible when based on a data-base of frozen free atomic electron densities, $\{n_I(\mathbf{r})\}$. The “Hartree” matrix is given by,

$$M_{IJ}^H = \begin{cases} 0 & \forall I = J, \\ \int d\mathbf{r} d\mathbf{r}' \frac{n_I(\mathbf{r})n_J(\mathbf{r}')}{r_{IJ}} & \forall I \neq J, \end{cases} \quad (22)$$

the “external” potential matrix is given by,

$$M_{IJ}^{ext} = \begin{cases} 0 & \forall I = J, \\ \int d\mathbf{r} \frac{n_I(\mathbf{r})Z_J}{r_{IJ}} & \forall I \neq J, \end{cases} \quad (23)$$

and the “kinetic” matrix is given within

$$M_{IJ}^{kin} = \begin{cases} C_F \int d\mathbf{r} n_I(\mathbf{r})^{5/3} & \forall I = J, \\ 0 & \forall I \neq J, \end{cases} \quad (24)$$

where C_F is a constant, and atomic integrals are evaluated over all of space. If need be, the kinetic matrix

could even be extended by the von Weizsäcker correction term, $\frac{1}{8} \int d\mathbf{r} |\nabla n_I(\mathbf{r})|^2 / n_I(\mathbf{r})^{136}$. Summation of all entries in the matrix and addition of the off-diagonal Coulomb-matrix entries would yield the corresponding exact DFT energy for frozen atomic electron densities. Unfortunately, preliminary training on atomization energies of the GDB-7 data set¹² indicates that neither use of the Lennard-Jones nor of the Thomas Fermi matrices leads to any significant improvement in predictive accuracy when compared to the original Coulomb-matrix representation in Ref.¹²¹. A possible explanation for this non-intuitive result is that these more sophisticated descriptors are no longer monotonic functions in geometries and stoichiometries—in contrast to the Coulomb matrix.

An alternative new descriptor, entirely consistent with the first principles view on CCS, has recently been proposed¹³⁷. Each atom I in the molecule is represented by its nuclear charge multiplied with a cosine term that contains a radial distribution function of atom I with respect to all other atoms J . Summing up the atomic contributions yields a Fourier series of radial distribution functions which, because of the superposition principle, is not only unique for each compound, but also invariant with respect to molecular rotations, translations, and atom indexing.

$$M(d) = \sum_I^{N_J} Z_J^n \cos\left[\frac{1}{Z_J} \sum_I^{N_I} Z_I e^{-(d-d_{IJ})^2/\sigma}\right], \quad (25)$$

where $d_{IJ} = |\mathbf{R}_I - \mathbf{R}_J|$, and n and σ are hyper parameters that can be optimized. This descriptor has units of charge ^{n} , d has units of distance and goes from zero beyond the largest interatomic distance. As in the case of the Coulomb-matrix described above, the environment of large or condensed systems can be accounted for by choosing N_I to be larger than N_J . The reader is referred to the original paper for further details¹³⁷.

VI. CONCLUDING REMARKS

We have reviewed a notion of chemical compound space that is consistent with any *ab initio* approach to

atomistic simulations. Starting from an energy hierarchy, variations in nuclear charge distributions have been discussed, followed by order-parameter based interpolation approaches, and statistical learning methods. The concepts presented offer a seamless and rigorous framework to unify electronic structure theory with rigorous rational as well as combinatorial compound design efforts. This view of chemical space is advantageous for several reasons, (i) equipped with such a notion, important fundamental questions can be tackled in the future, including rigorous definitions of diversity in CCS, property transferability, uncertainty, and selection bias in training sets; (ii) transferability and applicability typical for the black-box characteristics and the accuracy of *ab initio* calculations can be achieved; (iii) a mathematically, physically, and chemically rigorous notion of relevant input variables enables the application of sophisticated property optimization algorithms. Ultimately, efforts along these lines promise to lead to “the right compound for the right reason”, promising to replace by systematic engineering protocols the heuristics and serendipity on which most, if not all, of the past compound discoveries have relied.

VII. ACKNOWLEDGMENTS

The author is thankful for helpful discussions with C. Anderson, M. Cuendet, R. A. DiStasio, Jr., J. R. Hammond, F. Kiraly, A. Knoll, G. Montavon, J. E. Moussa, K. R. Müller, B. C. Rinderspacher, M. Rupp, A. Tkatchenko, D. Truhlar, M. Tuckerman, A. Vazquez-Mayagoitia. The many participants of the 2011-program “Navigating Chemical Compound Space for Materials and Bio Design” at the Institute for Pure and Applied Mathematics, UCLA, are also greatly acknowledged. This research used resources of the Argonne Leadership Computing Facility at Argonne National Laboratory, which is supported by the Office of Science of the U.S. DOE under contract DE-AC02-06CH11357.

* Electronic address: anatole@alcf.anl.gov

¹ P. Kirkpatrick and C. Ellis, *Nature* **432**, 823 (2004).

² K. Burke (2011), “Any method whose parametrization does not depend on the chemical system being studied can be called an *ab initio* method.” Oral communication, IPAM, UCLA.

³ T. Helgaker, P. Jørgensen, and J. Olsen, *Molecular Electronic-Structure Theory* (John Wiley & Sons, LTD, 2000).

⁴ M. E. Tuckerman, *Statistical mechanics: Theory and molecular simulation* (Oxford University Press, 2010).

⁵ P. Hohenberg and W. Kohn, *Phys. Rev.* **136**, B864 (1964).

⁶ A. D. Becke, *Phys. Rev. A* **38**, 3098 (1988).

⁷ C. Lee, W. Yang, and R. G. Parr, *Phys. Rev. B* **37**, 785 (1988).

⁸ H. R. Henze and C. M. Blair, *J. Am. Chem. Soc.* **53**, 3077 (1931).

⁹ D. Perry, *J. Am. Chem. Soc.* **54**, 2918 (1932).

¹⁰ L. Bytautas and D. Klein, *J. Chem. Inf. Comp. Sci.* **38**, 1063 (1998).

¹¹ J. Braun, R. Gugisch, A. Kerber, R. Laue, M. Meringer, and C. Rücker, *J. Chem. Inf. Comp. Sci.* **44**, 542 (2004).

¹² L. C. Blum and J.-L. Reymond, *J. Am. Chem. Soc.* **131**, 8732 (2009).

- ¹³ A. C. T. van Duin, S. Dasgupta, F. Lorant, , and W. A. Goddard III, *J. Phys. Chem. A* **105**, 9396 (2001).
- ¹⁴ A. K. Rappé, C. J. Casewit, K. S. Colwell, W. A. Goddard III, and W. M. Skid, *J. Am. Chem. Soc.* **114**, 10024 (1992).
- ¹⁵ P. Geerlings, F. D. Proft, and W. Langenaeker, *Chem. Rev.* **103**, 1793 (2003).
- ¹⁶ R. G. Parr and W. Yang, *Density functional theory of atoms and molecules* (Oxford Science Publications, 1989).
- ¹⁷ O. A. von Lilienfeld and M. E. Tuckerman, *J. Chem. Phys.* **125**, 154104 (2006).
- ¹⁸ J. F. Capitani, R. F. Nalewajski, and R. G. Parr, *J. Chem. Phys.* **76**, 568 (1982).
- ¹⁹ O. A. von Lilienfeld, R. Lins, and U. Rothlisberger, *Phys. Rev. Lett.* **95**, 153002 (2005).
- ²⁰ V. Marcon, O. A. von Lilienfeld, and D. Andrienko, *J. Chem. Phys.* **127**, 064305 (2007).
- ²¹ J. G. Kirkwood, *J. Chem. Phys.* **3**, 300 (1935).
- ²² O. A. von Lilienfeld and M. E. Tuckerman, *J. Chem. Theory Comput.* **3**, 1083 (2007).
- ²³ H. Hellmann, *J. Chem. Phys.* **3**, 61 (1935).
- ²⁴ H. Hellmann, *J. Chem. Phys.* **4**, 324 (1936).
- ²⁵ J. C. Phillips and L. Kleinman, *Phys. Rev.* **116**, 287 (1959).
- ²⁶ J. D. Weeks and S. A. Rice, *J. Chem. Phys.* **49**, 2741 (1968).
- ²⁷ G. B. Bachelet, D. R. Hamann, and M. Schluter, *Phys. Rev. B* **26**, 4199 (1982).
- ²⁸ P. A. Christiansen, Y. S. Lee, and K. S. Pitzer, *J. Chem. Phys.* **71**, 4445 (1979).
- ²⁹ P. Pulay, *Mol. Phys.* **229** (1969).
- ³⁰ C. Hartwigsen, S. Goedecker, and J. Hutter, *Phys. Rev. B* **58**, 3641 (1998).
- ³¹ M. M. Rieger and P. Vogl, *Phys. Rev. B* **52**, 16567 (1995).
- ³² B. Baumeier, P. Krüger, and J. Pollmann, *Phys. Rev. B* **73**, 195205 (2006).
- ³³ O. A. von Lilienfeld, I. Tavernelli, U. Rothlisberger, and D. Sebastiani, *J. Chem. Phys.* **122**, 014113 (2005).
- ³⁴ G. A. DiLabio, M. M. Hurley, and P. A. Christiansen, *J. Chem. Phys.* **116**, 9578 (2002).
- ³⁵ O. A. von Lilienfeld, I. Tavernelli, U. Rothlisberger, and D. Sebastiani, *Phys. Rev. Lett.* **93**, 153004 (2004).
- ³⁶ E. Torres and G. A. DiLabio, *J. Phys. Chem. Lett.* **3**, 1738 (2012).
- ³⁷ N. E. Christensen, *Phys. Rev. B* **30**, 5753 (1984).
- ³⁸ D. Segev and A. Janotti and C. G. Van de Walle, *Phys. Rev. B* **75**, 35201 (2007).
- ³⁹ K. Leung, S. B. Rempe, and O. A. von Lilienfeld, *J. Chem. Phys.* **130**, 204507 (2009).
- ⁴⁰ M. Sulpizi and M. Sprik, *Phys. Chem. Chem. Phys.* **10**, 5238 (2008).
- ⁴¹ D. Alfè, M. J. Gillan, and G. D. Price, *Nature* **405**, 172 (2000).
- ⁴² D. Sheppard, G. Henkelman, and O. A. von Lilienfeld, *J. Chem. Phys.* **133**, 084104 (2010).
- ⁴³ S. Goedecker, M. Teter, and J. Hutter, *Phys. Rev. B* **54**, 1703 (1996).
- ⁴⁴ M. Krack, *Theor. Chim. Acta* **114**, 145 (2005).
- ⁴⁵ F. Weigend, C. Schrod, and R. Ahlrichs, *J. Chem. Phys.* **121**, 10380 (2004).
- ⁴⁶ C. Cardenas, W. Tiznado, P. W. Ayers, and P. Fuentealba, *J. Phys. Chem. A* **115**, 2325 (2011).
- ⁴⁷ M. Lesiuk, R. Balawender, and J. Zachara, *J. Chem. Phys.* **136**, 034104 (2012).
- ⁴⁸ M. Wang, X. Hu, D. N. Beratan, and W. Yang, *J. Am. Chem. Soc.* **128**, 3228 (2006).
- ⁴⁹ D. Xiao, W. Yang, and D. N. Beratan, *J. Chem. Phys.* **129**, 044106 (2008).
- ⁵⁰ X. Hu, D. N. Beratan, and W. Yang, *J. Chem. Phys.* **129**, 064102 (2008).
- ⁵¹ D. Balamurugan, W. Yang, and D. N. Beratan, *J. Chem. Phys.* **129**, 174105 (2008).
- ⁵² S. Keinan, M. J. Therien, D. N. Beratan, and W. Yang, *J. Phys. Chem. A* **112**, 12203 (2008).
- ⁵³ B. C. Rinderspacher, J. Andzelm, A. Rawlett, J. Dougherty, D. N. Beratan, and W. Yang, *J. Chem. Theory Comput.* **5**, 3321 (2009).
- ⁵⁴ S. R. Marder, D. N. Beratan, and L.-T. Cheng, *Science* **252**, 103 (1991).
- ⁵⁵ C. Kuhn and D. N. Beratan, *J. Phys. Chem.* **100**, 10595 (1996).
- ⁵⁶ M. d'Avezac and A. Zunger, *Phys. Rev. B* **78**, 064102 (2008).
- ⁵⁷ N. Sablon, F. D. Proft, P. W. Ayers, and P. Geerlings, *J. Chem. Theory Comput.* **6**, 3671 (2010).
- ⁵⁸ W. Yang, A. J. Cohen, F. D. Proft, and P. Geerlings, *J. Chem. Phys.* **136**, 144110 (2012).
- ⁵⁹ J. P. Perdew, R. G. Parr, M. Levy, and J. L. Balduz, *Phys. Rev. Lett.* **49**, 1691 (1982).
- ⁶⁰ J. P. Perdew and M. Levy, *Phys. Rev. Lett.* **51**, 1884 (1983).
- ⁶¹ P. Mori-Sánchez, A. J. Cohen, and W. Yang, *Phys. Rev. Lett.* **102**, 066403 (2009).
- ⁶² L. A. Constantin, J. C. Snyder, J. P. Perdew, and K. Burke, *J. Chem. Phys.* **133**, 241103 (2010).
- ⁶³ Ralph G. Pearson, *J. Climate* **64**, 561 (1987).
- ⁶⁴ P. G. Mezey, *J. Am. Chem. Soc.* **107**, 3100 (1985).
- ⁶⁵ R. P. Feynman, *Phys. Rev.* **56**, 340 (1939).
- ⁶⁶ O. A. von Lilienfeld, *J. Chem. Phys.* **131**, 164102 (2009).
- ⁶⁷ P. E. Smith and W. F. van Gunsteren, *J. Chem. Phys.* **100**, 577 (1994).
- ⁶⁸ A. Putrino, D. Sebastiani, and M. Parrinello, *J. Chem. Phys.* **113**, 7102 (2000).
- ⁶⁹ A. Beste, R. J. Harrison, and T. Yanai, *J. Phys. Chem.* **125**, 074101 (2006).
- ⁷⁰ A. Pérez and O. A. von Lilienfeld, *J. Chem. Theory Comput.* **7**, 2358 (2011).
- ⁷¹ I.-C. Lin, M. D. Coutinho-Neto, C. Felsenheimer, O. A. von Lilienfeld, I. Tavernelli, and U. Rothlisberger, *Phys. Rev. B* **75**, 205131 (2007).
- ⁷² M. Stiborová, J. Sejbál, L. Borek-Dohalska, D. Aimova, J. Poljakova, K. Forsterova, M. Rupertova, J. Wiesner, J. Hudecek, M. Wiessler, et al., *Cancer Res.* **64**, 8374 (2004).
- ⁷³ D. Mousset, R. Rabot, P. Bouyssou, G. Coudert, and I. Gillaizeau, *Tetrahedron Lett.* **51**, 3987 (2010).
- ⁷⁴ A. Canals, M. Purciolas, J. Aymami, and M. Coll, *Acta Cryst.* **D61**, 1009 (2005).
- ⁷⁵ C. Bissantz, B. Kuhn, and M. Stahl, *J. Medic. Chem.* **53**, 5061 (2010).
- ⁷⁶ M. Kolar, T. Kubar, and P. Hobza, *J. Phys. Chem. B* **114**, 13446 (2010).
- ⁷⁷ B. Tidor, *J. Phys. Chem.* **97**, 1069 (1993).
- ⁷⁸ C. Oostenbrink and W. F. van Gunsteren, *Proc. Natl. Acad. Sci. USA* **102**, 6750 (2005).
- ⁷⁹ C. Oostenbrink, *J. Comp. Chem.* **30**, 212 (2009).
- ⁸⁰ R. Iftimie, P. Minary, and M. E. Tuckerman, *Proc. Natl. Acad. Sci. USA* **102**, 6654 (2005).

- ⁸¹ A. Laio, J. VandeVondele, and U. Rothlisberger, *J. Chem. Phys.* **116**, 6941 (2002).
- ⁸² I.-C. Lin, O. A. von Lilienfeld, M. D. Coutinho-Neto, I. Tavernelli, and U. Rothlisberger, *J. Phys. Chem. B* **111**, 14346 (2007).
- ⁸³ O. A. von Lilienfeld and A. Tkatchenko, *J. Chem. Phys.* **132**, 234109 (2010).
- ⁸⁴ R. A. DiStasio, O. A. von Lilienfeld, and A. Tkatchenko, *Proc. Natl. Acad. Sci. USA* **109**, 14791 (2012).
- ⁸⁵ Prize (2011), A prize award of the equivalent of an ounce of gold was announced during the Navigating Chemical Compound Space program in spring 2011 at the Institute of Pure and Applied Mathematics, UCLA. The prize is for finding an interpolating transform of two iso-electronic Hamiltonians such that the potential energy becomes linear in the interpolating order parameter. The ounce of gold is currently held in the form of 100 shares of iShares Trust fund (NYSEARCA:IAU), and will be dispensed in cash at instantaneous exchange rate upon recognition of a valid solution by a prize-committee. Apart from the author, the prize-committee consists of Profs. K. Burke, G. Henkelman, K. R. Müller, and M. E. Tuckerman. Contact the author regarding donations to increase the prize award. For more information, see <http://www.alcf.anl.gov/~anatole>.
- ⁸⁶ L. P. Hammett, *Chem. Rev.* **17**, 125 (1935).
- ⁸⁷ L. P. Hammett, *J. Am. Chem. Soc.* **59**, 96 (1937).
- ⁸⁸ D. G. Pettifor, *J. Phys. C: Solid State Phys.* **19**, 285 (1986).
- ⁸⁹ J. F. Gonthier, S. N. Steinmann, M. D. Woodrich, and C. Corminboeuf, *Chem. Soc. Rev.* **41**, 4671 (2012).
- ⁹⁰ T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: data mining, inference and prediction*, Springer series in statistics (Springer, New York, N.Y., 2001).
- ⁹¹ B. Schölkopf and A. J. Smola, *Learning with Kernels* (MIT Press, Cambridge, 2002).
- ⁹² T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning. Data Mining, Inference, and Prediction* (Springer, New York, 2009), 2nd ed.
- ⁹³ K.-R. Müller, S. Mika, G. Rätsch, K. Tsuda, and B. Schölkopf, *IEEE Transactions on Neural Networks* **12**, 181 (2001).
- ⁹⁴ G. Csányi, T. Albaret, M. C. Payne, and A. D. Vita, *Phys. Rev. Lett.* **93**, 175503 (2004).
- ⁹⁵ P. Maurer, A. Laio, H. W. Hugosson, M. C. Colombo, and U. Rothlisberger, *J. Chem. Theory Comput.* **3**, 628 (2007).
- ⁹⁶ L. P. Wang and T. V. Voorhis, *J. Chem. Phys.* **133**, 231101 (2010).
- ⁹⁷ L. Hu, X. Wang, L. Wong, and G. Chen, *J. Chem. Phys.* **119**, 11501 (2003).
- ⁹⁸ X. Zheng, L. Hu, X. Wang, and G. Chen, *Chem. Phys. Lett.* **390**, 186 (2004).
- ⁹⁹ A. Brown, B. J. Braams, K. Christoffel, Z. Jin, and J. M. Bowman, *J. Chem. Phys.* **119**, 8790 (2003).
- ¹⁰⁰ S. Lorenz, A. Gross, and M. Scheffler, *Chem. Phys. Lett.* **395**, 210 (2004).
- ¹⁰¹ J. Behler and M. Parrinello, *Phys. Rev. Lett.* **98**, 146401 (2007).
- ¹⁰² C. M. Handley and P. L. A. Popelier, *J. Chem. Theory Comput.* **5**, 1474 (2009).
- ¹⁰³ J. Behler, R. Martonak, D. Donadio, and M. Parrinello, *Phys. Rev. Lett.* **100**, 185501 (2008).
- ¹⁰⁴ J. Behler, *Phys. Chem. Chem. Phys.* **13**, 17930 (2011).
- ¹⁰⁵ R. M. Balabin and E. I. Lomakina, *Phys. Chem. Chem. Phys.* **13**, 11710 (2011).
- ¹⁰⁶ A. P. Bartók, M. C. Payne, R. Kondor, and G. Csányi, *Phys. Rev. Lett.* **104**, 136403 (2010).
- ¹⁰⁷ S. Curtarolo, D. Morgan, K. Persson, J. Rodgers, and G. Ceder, *Phys. Rev. Lett.* **91**, 135503 (2003).
- ¹⁰⁸ W. Setyawan and S. Curtarolo, *Comp. Mat. Sci.* **49**, 299 (2010).
- ¹⁰⁹ G. Hautier, C. C. Fischer, A. Jain, T. Mueller, and G. Ceder, *Chem. Mater.* **22**, 3762 (2010).
- ¹¹⁰ G. R. Hutchison, M. A. Ratner, and T. J. Marks, *J. Am. Chem. Soc.* **127**, 2339 (2005).
- ¹¹¹ M. Misra, D. Andrienko, B. Baumeier, J.-L. Faulon, and O. A. von Lilienfeld, *J. Chem. Theory Comput.* **7**, 2549 (2011).
- ¹¹² J. C. Snyder, M. Rupp, K. Hansen, K.-R. Müller, and K. Burke, *Phys. Rev. Lett.* **108**, 253002 (2012).
- ¹¹³ Z. D. Pozun, K. Hansen, D. Sheppard, M. Rupp, K.-R. Müller, and G. Henkelman, *J. Chem. Phys.* **136**, 174101 (2012).
- ¹¹⁴ J. Wellendorff, K. T. Lundgaard, A. Møgelhøj, V. Petzold, D. D. Landis, J. K. Nørskov, T. Bligaard, and K. W. Jacobsen, *Phys. Rev. B* **85**, 235149 (2012).
- ¹¹⁵ G. Schneider, *Nature Reviews* **9**, 273 (2010).
- ¹¹⁶ O. Ivanciuc, *J. Chem. Inf. Comp. Sci.* **40**, 1412 (2000).
- ¹¹⁷ R. Todeschini and V. Consonni, *Handbook of Molecular Descriptors* (Wiley-VCH, Weinheim, 2009).
- ¹¹⁸ J. Braun, A. Kerber, M. Meringer, and C. Rücker, *MATCH* **54**, 163 (2005).
- ¹¹⁹ J.-L. Faulon, D. P. Visco, Jr., and R. S. Pophale, *J. Chem. Inf. Comp. Sci.* **43**, 707 (2003).
- ¹²⁰ X. Feng, V. Marcon, W. Pisula, M. R. Hansen, J. Kirkpatrick, F. Grozema, D. Andrienko, K. Kremer, and K. Mullen, *Nature Materials* **8**, 421 (2009).
- ¹²¹ M. Rupp, A. Tkatchenko, K.-R. Müller, and O. A. von Lilienfeld, *Phys. Rev. Lett.* **108**, 058301 (2012).
- ¹²² W. Kohn and L. J. Sham, *Phys. Rev.* **140**, A1133 (1965).
- ¹²³ A. D. Becke, *J. Chem. Phys.* **98**, 5648 (1993).
- ¹²⁴ J. P. Perdew, M. Ernzerhof, and K. Burke, *J. Chem. Phys.* **105**, 9982 (1996).
- ¹²⁵ M. Ernzerhof and G. E. Scuseria, *J. Chem. Phys.* **110**, 5029 (1999).
- ¹²⁶ B. J. Lynch and D. G. Truhlar, *J. Phys. Chem. A* **107**, 3898 (2003).
- ¹²⁷ G. Montavon, K. Hansen, S. Fazli, M. Rupp, F. Biegler, A. Ziehe, A. Tkatchenko, O. A. von Lilienfeld, and K.-R. Müller, NIPS proceedings (2013), accepted.
- ¹²⁸ G. Montavon, M. Rupp, V. Gobre, A. Vazquez-Mayagoitia, K. Hansen, A. Tkatchenko, K.-R. Müller, and O. A. von Lilienfeld (2012), submitted.
- ¹²⁹ Using 2 as powers of Z (the energy of the hydrogenic atom) did not change performance.
- ¹³⁰ A. L. Patterson, *Nature* **143**, 939 (1939).
- ¹³¹ J. E. Moussa, *Phys. Rev. Lett.* **109**, 059801 (2012).
- ¹³² M. Rupp, A. Tkatchenko, K.-R. Müller, and O. A. von Lilienfeld, *Phys. Rev. Lett.* **109**, 059802 (2012).
- ¹³³ D. C. Ciresan, U. Meier, L. M. Gambardella, and J. Schmidhuber, *Neural Computation* **22**, 3207 (2010).
- ¹³⁴ J. Behler, *J. Chem. Phys.* **134**, 074106 (2011).
- ¹³⁵ S. Doraiswamy, J. Bender, G. V. Candler, Y. Paukku, K. Yang, Z. Varga, and D. G. Truhlar (2012), to be published.
- ¹³⁶ H. Chen and A. Zhou, *Numer. Math. Theor. Meth. Appl.*

1, 1 (2008).

¹³⁷ O. A. von Lilienfeld and A. Knoll (2012), submitted.

¹³⁸ Include the mass of the nuclei as an additional variable if also dynamical properties and nuclear quantum effects

are to be accounted for.

¹³⁹ We assume the reduced mass to equate the mass of the electron